

هوش مصنوعی همه ما را شبیه به هم می‌کند!

شواهد نوظهور نشان می‌دهد که خروجی مدل‌های زبانی بزرگ می‌تواند متن و حتی افکار کاربران انسانی را شکل دهد.



شما هر روز به نظر می‌دهید که خود را از بقیه متمایز می‌دانید. اما در واقع، افکار کاربران انسانی به شکل قابل توجهی در حال همگونی است. در قلب مدل‌های هوش مصنوعی امروزی، مجموعه‌های عظیمی از داده‌های آموزشی قرار دارد. متن‌ها، ویدیوها و تصاویر ساخته شده توسط انسان‌ها که برای آموزش مدل‌ها به کار می‌روند تا الگوها را شناسایی کرده و محتوا تولید کنند. بدون شک انسان‌ها در حال آموزش سیستم‌های هوش مصنوعی هستند، اما آیا هوش مصنوعی نیز در حال آموزش دادن ماست؟

به نقل از نیچر، تعداد رو به رشدی از مقالات گزارش می‌دهند که افراد تمایل دارند الگوهای نوشتاری، روش‌های استدلال و حتی دیدگاه‌ها را از مدل‌های زبانی بزرگی (LLM) که استفاده می‌کنند، دریافت کنند. برخی پژوهشگران می‌گویند این تأثیر می‌تواند به نوعی یکنواختی در نوشتار انسانی منجر شود و هشدار می‌دهند که این اثر حتی ممکن است به متونی که توسط افرادی نوشته می‌شود که مستقیماً از هوش مصنوعی استفاده نمی‌کنند نیز سرایت کند.

ژویار سوراتی، دانشمند علوم رایانه در دانشگاه کالیفرنیا جنوبی در لس‌آنجلس و یکی از نویسندگان مقاله‌ای در این مورد استدلال می‌کند که مدل‌های زبانی بزرگ در حال همگن‌سازی گفتمان انسانی هستند و می‌گویند: اگر اطرافیان شما با این مدل‌ها تعامل داشته باشند و سبک نوشتاری، دیدگاه‌ها و شیوه‌های استدلال آن‌ها را بپذیرند، در نهایت این الگوها آن قدر شما را احاطه می‌کنند که به نظر می‌رسد این روش، شکل اجتماعی درست برای بیان اطلاعات است. با این حال، برخی دیگر معتقدند ذهن انسان ممکن است همچنان در برابر این اثر هموارکننده هوش مصنوعی مقاومت کند. در یک مطالعه که در ماه نوامبر به صورت پیش‌چاپ در سرور arXiv منتشر شد، نویسندگان گروه‌هایی از نویسندگان را شناسایی کردند که «نشانه‌های سبکی متمایز و انسانی» خود را حفظ می‌کنند و احتمالاً اصالت را به سود کارایی ارائه شده توسط هوش مصنوعی ترجیح می‌دهند. این مطالعه هنوز داوری هم‌تا نشده است.

بررسی عمیق‌تر موضوع

در یک پیش‌چاپ دیگر که سال گذشته در arXiv منتشر شد و هنوز داوری نشده، سوراتی و همکارانش پست‌های ردیت، محتوای خبری و مطالعات پیش‌چاپ را پیش و پس از راه‌اندازی چت‌جی‌پی‌تی در سال ۲۰۲۲ تحلیل کردند. گروه پژوهشی دریافت که متن‌های منتشر شده پس از عرضه این پلتفرم، از نظر سبک نوشتاری تنوع کمتری نسبت به متن‌های قبلی دارند. در مقاله‌ای جدید، نویسندگان استدلال می‌کنند که این پدیده بر دیدگاه‌ها و شیوه‌های استدلال افراد نیز تأثیر می‌گذارد. آن‌ها به یک پیش‌چاپ بدون داوری در سال ۲۰۲۳ اشاره می‌کنند که در آن شرکت‌کنندگان با مدل‌های زبانی‌ای تعامل داشتند که دیدگاه‌های مثبت یا منفی درباره شبکه‌های اجتماعی بیان می‌کردند. پس از این مواجهه، دیدگاه‌های خود شرکت‌کنندگان به سمت دیدگاه‌های تولیدشده توسط مدل‌ها تغییر کرد.

اولیور هاوزر، پژوهشگر اقتصاد و هوش مصنوعی در دانشگاه اکستر بریتانیا، می‌گوید نویسندگان نکته قابل قبولی مطرح می‌کنند که افراد می‌توانند از هوش مصنوعی سود ببرند و آن این است که به شما کمک می‌کند بهتر بنویسید و برای دیگران قابل فهم‌تر باشید. اما او هشدار می‌دهد: به محض اینکه این پذیرش فراگیر شود، این جمع است که بیشترین آسیب را می‌بیند.

ورود به عرصه سیاست

در مطالعه‌ای که به تازگی در مجله Science Advances منتشر شده، پژوهشگران دریافتند که دیدگاه‌های افراد درباره مسائل اجتماعی شروع به بازتاب دیدگاه‌هایی می‌کند که از یک ابزار هوش مصنوعی دریافت کرده‌اند. شرکت‌کنندگان از دستیارهای هوش مصنوعی برای نوشتن درباره موضوعات اجتماعی-سیاسی مانند مجازات اعدام استفاده کردند. پس از آن، نگرش‌های آن‌ها بیشتر به آنچه مدل‌های زبانی نوشته بودند در مقایسه با گروه کنترل که از هوش مصنوعی استفاده نکرده بودند، شبیه شد. استرلینگ ویلیامز-سسی، یکی از نویسندگان این مطالعه و دانشمند اطلاعات در دانشگاه کرنل در ایتاکا، نیویورک، می‌گوید این اثر می‌تواند در نهایت تنوع دیدگاه‌های سیاسی را نیز کاهش دهد. البته میزان این تأثیر به گرایش‌هایی بستگی دارد که مدل‌های مختلف زبانی از خود نشان می‌دهند.

نکته مهم این است که شرکت‌کنندگان متوجه نشدند که تحت تأثیر چت‌بات‌ها قرار گرفته‌اند. حتی زمانی که به آن‌ها گفته شد هوش مصنوعی ممکن است دیدگاه‌هایشان را سوگیرانه کند، نتایج تغییری نکرد. ویلیامز-سسی می‌گوید: در حال حاضر نمی‌دانیم چگونه می‌توان از این موضوع جلوگیری کرد. و احتمالاً راه حل به سادگی ارائه یک هشدار ساده به کاربران نیست. هاوزر همچنین می‌گوید یکنواختی تحمیل‌شده توسط ابزارهای هوش مصنوعی می‌تواند تفکر علمی را نیز محدود کند: ممکن است یک ایده دیوانه‌وار را از دست بدهیم که در ابتدا غیرمنطقی به نظر می‌رسد، اما در نهایت همان چیزی است که برای یک جهش علمی به آن نیاز داریم. مقاله‌ای که در ژانویه در نیچر منتشر شد نشان داد دانشمندانی که از ابزارهای هوش مصنوعی در پژوهش‌های خود استفاده می‌کنند، نسبت به کسانی که استفاده نمی‌کنند، تمایل دارند روی مجموعه محدودتری از حوزه‌ها تمرکز کنند.

سبک شخصی

با این حال، همه مطالعات نشان نمی‌دهند که استفاده از هوش مصنوعی به یکنواختی منجر می‌شود. مقاله‌ای در ماه نوامبر گزارش داد که نویسندگان انسانی هنگام تعامل با هوش مصنوعی به شیوه‌های متفاوتی تکامل پیدا می‌کنند. نوشتار برخی افراد واقعاً به مدل‌های هوش مصنوعی شبیه‌تر می‌شود؛ در حالی که دیگران سبک شخصی خود را حفظ می‌کنند یا حتی سبکی ایجاد می‌کنند که بیش از پیش از هوش مصنوعی متمایز است.

مطالعاتی که نگرانی‌هایی درباره از بین رفتن تنوع در نوشتار انسانی مطرح کرده‌اند، باعث تردید یون وان، پژوهشگر همکاری انسان و هوش مصنوعی در دانشگاه هیوستون-داون تاون تگزاس، شد. او می‌گوید: ما فکر کردیم شاید این طور نباشد. وان و همکارانش ۱۰ شخصیت هوش مصنوعی منحصر به فرد با پیشینه‌های فرهنگی و سبک‌های فکری متفاوت ایجاد کردند که طرح داستان‌هایی تولید می‌کردند و شرکت‌کنندگان انسانی بر اساس آن‌ها داستان‌های خلاقانه می‌نوشتند. تحلیل نویسندگان که سال گذشته به صورت پیش‌چاپ منتشر شد و هنوز داوری نشده، نشان داد که داستان‌های حاصل، سطحی از تنوع مشابه داستان‌هایی را داشتند که بدون کمک هوش مصنوعی نوشته شده بودند. وان می‌گوید به دلیل همین پاسخ‌های متنوع، این شخصیت‌های سفارشی می‌توانند راهی برای کاهش یکنواختی باشند.